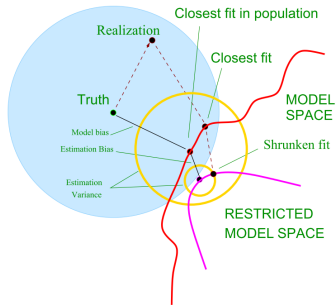


Introduction to Machine Learning

Regularization

Bias-variance Tradeoff



Learning goals

- Understand the bias-variance trade-off
- Know the definition of model bias, estimation bias, and estimation variance

BIAS-VARIANCE TRADEOFF

In this slide set, we will visualize the bias-variance trade-off.

First, we start with a DGP \mathbb{P}_{xy} and a suitable loss function $L : \mathbb{R}^g \times \mathbb{R}^g \rightarrow \mathbb{R}$ where \mathbb{R}^g is numerical encoding of \mathcal{Y} . We measure the distance between models $f : \mathcal{X} \rightarrow \mathbb{R}^g$ via

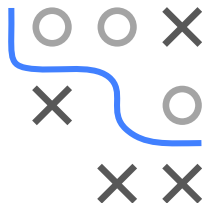
$$d(f, f') = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathbf{x}}} [L(f(\mathbf{x}), f'(\mathbf{x}))].$$

We restrict our attention to losses for which d becomes a metric, e.g., L1-loss, L2-loss, etc.

We define f_{true} as the risk minimizer such that

$$f_{\text{true}} \in \arg \min_{f \in \mathcal{H}_0} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(y, f(\mathbf{x}))]$$

where $\mathcal{H}_0 = \{f : \mathcal{X} \rightarrow \mathbb{R}^g \mid d(\underline{0}, f) < \infty\}$ and $\underline{0} : \mathcal{X} \rightarrow \{0\}$.



BIAS-VARIANCE TRADEOFF / 2

In practice, our model space \mathcal{H} usually is a proper subset of \mathcal{H}_0 and in general $f_{\text{true}} \notin \mathcal{H}$.

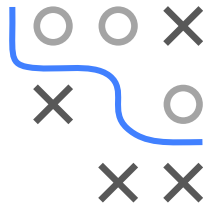
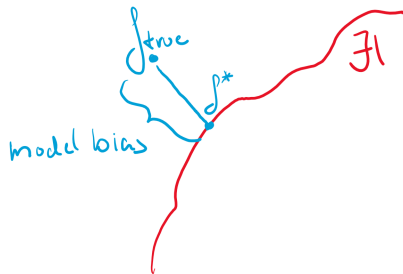
We define f^* as the risk minimizer in \mathcal{H} , i.e.,

$$f^* \in \arg \min_{f \in \mathcal{H}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(f(\mathbf{x}, y))].$$

It is the function in \mathcal{H} closest to f_{true} , and we call $d(f_{\text{true}}, f^*)$ the model bias.



BIAS-VARIANCE TRADEOFF / 3



BIAS-VARIANCE TRADEOFF / 4

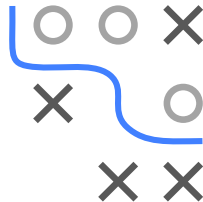
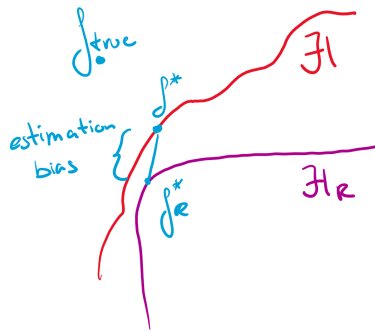
We can further restrict the model space such that \mathcal{H}_R is a proper subset of \mathcal{H} . We define f_R^* as the risk minimizer in \mathcal{H}_R , i.e.,

$$f_R^* \in \arg \min_{f \in \mathcal{H}_R} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(f(\mathbf{x}), y)] .$$

It is the function in \mathcal{H}_R closest to f_{true} , and we call $d(f_R^*, f^*)$ the estimation bias.



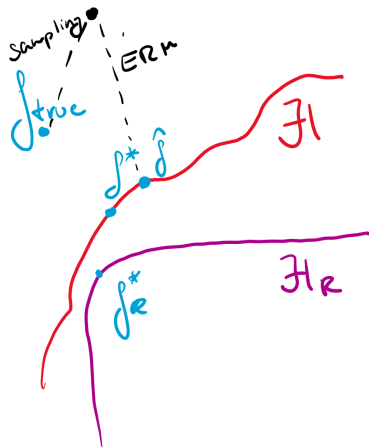
BIAS-VARIANCE TRADEOFF / 5



BIAS-VARIANCE TRADEOFF / 6

We sample a finite dataset $\mathcal{D} = (\mathbf{x}^{(i)}, y^{(i)})^n \in (\mathbb{P}_{xy})^n$ and find via ERM

$$\hat{f} \in \arg \min_{f \in \mathcal{H}} \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)})).$$



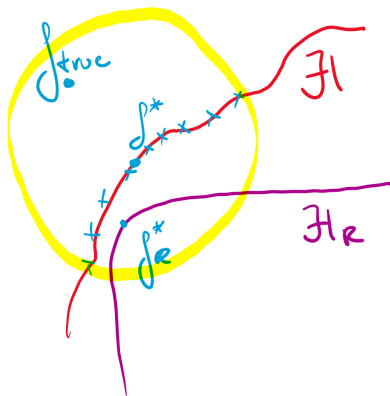
Note:

- $L : \mathcal{Y} \times \mathbb{R}^g \rightarrow \mathbb{R}$ is overloaded.
- The samples are only shown in the visualization for didactic purposes but are not an element of \mathcal{H} .



BIAS-VARIANCE TRADEOFF / 7

Let's assume that \hat{f} is an unbiased estimate of f^* (e.g., valid for linear regression), and we repeat the sampling process of \hat{f} .



- We can measure the spread of sampled \hat{f} around f^* via $\delta = \text{Var}_{\mathcal{D}} [d(f^*, \hat{f})]$ which we call the estimation variance.
- We visualize this as a circle around f^* with radius δ .



A 3x3 grid with a blue path starting at the top-left cell (0,0) and ending at the bottom-right cell (2,2). The path is composed of three segments: a horizontal segment from (0,0) to (1,0), a vertical segment from (1,0) to (1,1), and a diagonal segment from (1,1) to (2,2). The cells (0,1), (0,2), (1,2), and (2,0) are empty. The cells (1,0), (1,1), and (2,1) contain a black 'X'. The cells (0,0), (0,1), (0,2), (1,2), and (2,0) contain a grey circle.

Sampling

f

f^*

Error in f_k

H_k

$$\hat{f}_R \in \arg \min_{f \in \mathcal{H}_R} \sum_{i=1}^n L\left(y^{(i)}, \hat{f}(\mathbf{x}^{(i)})\right).$$

- We can measure the spread of sampled \hat{f}_R around f_R^* via $\delta = \text{Var}_{\mathcal{D}} \left[d(f^*, \hat{f}_R) \right]$ which we also call estimation variance.
- We observe that the increased bias results in a smaller estimation variance in \mathcal{H}_R compared to \mathcal{H} .